

# Accurate and Efficient Halo-based Galaxy Clustering Modelling with Simulations

Zheng Zheng<sup>1\*</sup> and Hong Guo<sup>2,1</sup>

<sup>1</sup> *Department of Physics and Astronomy, University of Utah, 115 South 1400 East, Salt Lake City, UT 84112, USA*

<sup>2</sup> *Shanghai Astronomical Observatory, Chinese Academy of Sciences, Shanghai 200030, China*

March 3, 2016

## ABSTRACT

Small- and intermediate-scale galaxy clustering can be used to establish the galaxy-halo connection to study galaxy formation and evolution and to tighten constraints on cosmological parameters. With the increasing precision of galaxy clustering measurements from ongoing and forthcoming large galaxy surveys, accurate models are required to interpret the data and extract relevant information. We introduce a method based on high-resolution  $N$ -body simulations to accurately and efficiently model the galaxy two-point correlation functions (2PCFs) in projected and redshift spaces. The basic idea is to tabulate all information of haloes in the simulations necessary for computing the galaxy 2PCFs within the framework of halo occupation distribution or conditional luminosity function. It is equivalent to populating galaxies to dark matter haloes and using the mock 2PCF measurements as the model predictions. Besides the accurate 2PCF calculations, the method is also fast and therefore enables an efficient exploration of the parameter space. As an example of the method, we decompose the redshift-space galaxy 2PCF into different components based on the type of galaxy pairs and show the redshift-space distortion effect in each component. The generalizations and limitations of the method are discussed.

**Key words:** cosmology: observations – cosmology: theory – galaxies: clustering – galaxies: distances and redshifts – galaxies: haloes – galaxies: statistics – large-scale structure of Universe

## 1 INTRODUCTION

Over the past two decades, large galaxy redshift surveys, such as the Sloan Digital Sky Survey (SDSS; York et al. 2000), the Two-Degree Field Galaxy Redshift Survey (2dFGRS; Colless 1999), the SDSS-III (Eisenstein et al. 2011), and the WiggleZ Dark Energy Survey (Blake et al. 2011), have enabled us to study in detail the large-scale structure of the universe probed by galaxies. Galaxy clustering has become a powerful tool to study galaxy formation and evolution and to learn about cosmology. An informative way to interpret galaxy clustering is to link galaxies to the underlying dark matter halo population, whose formation and evolution are dominated by gravitational interaction and whose properties are well understood with analytic models and  $N$ -body simulations.

The commonly adopted descriptions of the connection between galaxies and dark matter haloes include the halo occupation distribution (HOD; e.g. Jing et al. 1998; Peacock & Smith 2000; Seljak 2000; Scoccimarro et al. 2001; Berlind & Weinberg 2002; Berlind et al. 2003; Zheng et al. 2005) and the conditional luminosity function (CLF; e.g. Yang et al. 2003). The former specifies the

probability distribution of the number of galaxies in a given sample as a function of halo mass, together with the spatial and velocity distribution of galaxies inside haloes. The latter specifies the luminosity distribution of galaxies as a function of halo mass. Given a set of HOD or CLF parameters, with the halo population for a given cosmological model, galaxy clustering statistics can be predicted. Such frameworks have been successfully applied to galaxy clustering data to infer the connection between galaxy properties and halo mass (see e.g. van den Bosch et al. 2003; Zehavi et al. 2005, 2011; Zheng et al. 2007, 2009; Guo et al. 2014; Skibba et al. 2015) and to constrain cosmology (e.g. van den Bosch et al. 2003; Tinker et al. 2005; Cacciato et al. 2013; Reid et al. 2014). In particular, the main clustering statistic used is the two-point correlation function (2PCF) of galaxies, which is the focus of this paper as well.

Halo properties, like their mass function and spatial clustering (bias), can be understood analytically (e.g. Press & Schechter 1974; Mo et al. 1996; Sheth & Tormen 1999), and  $N$ -body simulations also enable accurate fitting formulae to be obtained (e.g. Jenkins et al. 2001; Tinker et al. 2008, 2010). Based on these, analytic models of galaxy 2PCF can be developed. The basic idea is to decompose the 2PCF into contributions from intra-halo and inter-

\* E-mail: zhengzheng@astro.utah.edu

halo galaxy pairs. The intra-halo component, or the one-halo term, represents the highly nonlinear part of the 2PCF. The inter-halo component, or the two-halo term, can be largely modelled by linear theory. Such analytic models have the advantage of being computationally inexpensive, and they can be used to efficiently probe the HOD/CLF and cosmology parameter space. However, as the precision of the 2PCF measurements in large galaxy surveys continues to improve, the requirement on the accuracy of the analytic models becomes more and more demanding. As pointed out in [Zheng \(2004a\)](#), an accurate model of the galaxy 2PCF needs to incorporate the nonlinear growth of the matter power spectrum (e.g. [Smith et al. 2003](#)), the halo exclusion effect, and the scale-dependent halo bias. In addition, the non-spherical shape of haloes should also be accounted for (e.g. [Tinker et al. 2005](#); [van den Bosch et al. 2013](#)). These are just factors to be taken into account in computing the real-space or projected 2PCFs. For redshift-space 2PCFs, more factors come into play. An accurate analytical description of the velocity field of dark matter haloes in the nonlinear or weakly nonlinear regime proves to be difficult and complex (e.g. [Tinker 2007](#); [Reid & White 2011](#); [Zu & Weinberg 2013](#)). Therefore, an accurate analytic model of redshift-space 2PCFs on small and intermediate scales is still not within reach.

The above complications faced by analytic models can all be avoided or greatly reduced if the 2PCF calculation is directly done with the outputs of  $N$ -body simulations. With the simulation, dark matter haloes can be identified, and their properties (mass, velocity, etc) can be obtained. For a given set of HOD/CLF parameters, one can populate haloes with galaxies accordingly (e.g. using dark matter particles as tracers) and form a mock galaxy catalog. The 2PCFs measured from the mock catalog are then the model predictions used to model the measurements from observations. Such a method of directly populating simulations have been developed and applied to model galaxy clustering data (e.g. [White et al. 2011](#); [Parejko et al. 2013](#)). This simulation-based model is attractive, as more and more large high-resolution  $N$ -body simulations emerge. It is also straightforward to implement. Once the mock catalog is produced, measuring the 2PCFs can be made fast (e.g. with tree code). However, populating haloes with a given set of HOD/CLF parameters is probably the most time-consuming step, as one needs to loop over all haloes of interest. In addition, information of individual haloes and tracer particles is needed, like their positions and velocities. Even with only a subset of all the particles in a high-resolution simulation, the amount of data can still be substantial.

The purpose of this paper is to introduce a method that takes the advantage of the simulation-based model, but being much more efficient in modelling galaxy clustering. The main idea is to decompose the galaxy 2PCFs and compress the information in the simulation by tabulating relevant clustering-related quantities of dark matter haloes. We also apply a similar idea to extend the commonly used sub-halo abundance matching method (SHAM; e.g. [Conroy et al. 2006](#)).

The paper is structured as follows. In Section 2, we formulate the method, within the HOD/CLF-like framework and within the halo/sub-halo framework. In Section 3, we show an example of modelling redshift-space 2PCFs, which also provides an understanding of the three-dimensional (3D) small- and intermediate-scale galaxy redshift-space 2PCF and its multipoles by decomposing them into the various components. In Section 4, we summarize the method and discuss possible generalizations and limitations.

## 2 SIMULATION-BASED METHOD OF CALCULATING GALAXY 2PCFS

In our simulation-based method, we divide haloes identified in  $N$ -body simulations into narrow bins of a given property, which determines galaxy occupancy. In the commonly used HOD/CLF, the property is the halo mass. In our presentation, we use halo mass as the halo variable, but the method can be generalized to any set of halo properties.

The basic idea of the method is to decompose the galaxy 2PCF into contributions from haloes of different masses, from one-halo and two-halo terms, and from different types of galaxy pairs (e.g. central-central, central-satellite, and satellite-satellite pairs). The decomposition also allows the separation between the halo occupation and halo clustering. The former relies on the specific HOD/CLF parameterization, while the latter can be calculated from the simulation. The method is to tabulate all relevant information about the latter for efficient calculation of galaxy 2PCFs and exploration of the HOD/CLF parameter space.

We first formulate the method in the HOD/CLF framework. We then apply the similar idea to the SHAM case, which provides a more general SHAM method.

### 2.1 Case with Simulation Particles

Let us start with a given  $N$ -body simulation and a given set of HOD/CLF parameters. To populate galaxies into a halo identified in the simulation, we can put one galaxy at the halo ‘centre’ as a central galaxy, according to the probability specified by the HOD/CLF parameters. Halo ‘centre’ should be defined to reflect galaxy formation physics. For example, a sensible choice is the position of potential minimum rather than centre of mass. For satellites, we can choose particles as tracers. In the usually adopted models, it is assumed that satellite galaxies follow dark matter particles inside haloes (e.g. [Zheng 2004a](#); [Tinker et al. 2005](#); [van den Bosch et al. 2013](#)), rooted in theoretical basis (e.g. [Nagai & Kravtsov 2005](#)). One can certainly modify the distribution profile as needed, and below we assume that the distribution of galaxies inside haloes has been specified and that the corresponding tracer particles have been selected for each halo.

We divide haloes in the simulation into  $N$  narrow mass bins and denote the mean number density of haloes in the mass bin  $\log M_i \pm d \log M_i / 2$  as  $\bar{n}_i$ . The mean number density of galaxies is computed as

$$\bar{n}_g = \sum_i \bar{n}_i [\langle N_{\text{cen}}(M_i) \rangle + \langle N_{\text{sat}}(M_i) \rangle], \quad (1)$$

where  $N_{\text{cen}}(M)$  and  $N_{\text{sat}}(M)$  are the occupation numbers of central and satellite galaxies in a halo of mass  $M$ ,  $\langle \rangle$  denotes the average over all haloes of this mass, and  $i = 1, \dots, N$ .

In the halo-based model, galaxy 2PCF  $\xi_{\text{gg}}$  is computed as the combination of two terms,  $\xi_{\text{gg}} = 1 + \xi_{\text{gg}}^{\text{1h}} + \xi_{\text{gg}}^{\text{2h}}$  ([Zheng 2004a](#)), where the one-halo term  $\xi_{\text{gg}}^{\text{1h}}$  (two-halo term  $\xi_{\text{gg}}^{\text{2h}}$ ) are from contributions of intra-halo (inter-halo) galaxy pairs. Following [Berlind & Weinberg \(2002\)](#), the one-halo term can be computed based on

$$\frac{1}{2} \bar{n}_g (\bar{n}_g d^3 \mathbf{r}) [1 + \xi_{\text{gg}}^{\text{1h}}(\mathbf{r})] = \sum_i \bar{n}_i \langle N_{\text{pair}}(M_i) \rangle f(\mathbf{r}; M_i) d^3 \mathbf{r}. \quad (2)$$

The left-hand side (LHS) is the number density of one-halo pairs with separation in the range  $\mathbf{r} \pm d\mathbf{r}/2$  from the definition of 2PCF.

The right-hand side (RHS) is the same quantity from counting one-halo pairs in each halo and the summation is over all the halo mass bins. Here  $\langle N_{\text{pair}}(M) \rangle$  is the total mean number of galaxy pairs in haloes of mass  $M$ , and  $f(\mathbf{r}; M)$  is the probability distribution of pair separation in haloes of mass  $M$ , i.e.  $f(\mathbf{r}; M)d^3\mathbf{r}$  is the probability of finding pairs with separation in the range  $\mathbf{r} \pm d\mathbf{r}/2$  in haloes of  $M$ . By further decomposing pairs into central-satellite (cen-sat) and satellite-satellite (sat-sat) pairs, we reach the following expression,

$$1 + \xi_{\text{gg}}^{\text{1h}}(\mathbf{r}) = \sum_i 2 \frac{\bar{n}_i}{\bar{n}_g^2} \langle N_{\text{cen}}(M_i) N_{\text{sat}}(M_i) \rangle f_{\text{cs}}(\mathbf{r}; M_i) + \sum_i \frac{\bar{n}_i}{\bar{n}_g^2} \langle N_{\text{sat}}(M_i) [N_{\text{sat}}(M_i) - 1] \rangle f_{\text{ss}}(\mathbf{r}; M_i). \quad (3)$$

The functions  $f_{\text{cs}}(\mathbf{r}; M)$  and  $f_{\text{ss}}(\mathbf{r}; M)$  are the probability distributions of one-halo cen-sat and sat-sat galaxy pair separation in haloes of mass  $M$ . They are normalized such that

$$\int f_{\text{cs}}(\mathbf{r}; M) d^3\mathbf{r} = 1 \quad \text{and} \quad \int f_{\text{ss}}(\mathbf{r}; M) d^3\mathbf{r} = 1. \quad (4)$$

Note that here and in what follows, the 2PCF can be either real-space, projected-space, redshift-space, or it can be the multipoles of the redshift-space 2PCF. The variable  $\mathbf{r}$  should be understood as pair separation in the corresponding space. For redshift-space clustering, we discuss how to specify velocity distribution of galaxies later.

To compute the two-halo term, we add up all possible two-halo galaxy pairs, following the 2PCF decomposition from different pair counts in [Zu et al. \(2008\)](#). Similar to equation (2), the total number density of two-halo pairs with separation in the range  $\mathbf{r} \pm d\mathbf{r}/2$  is

$$n_{\text{pair}, 2\text{h}} = \frac{1}{2} \bar{n}_g (\bar{n}_g d^3\mathbf{r}) [1 + \xi_{\text{gg}}^{\text{2h}}(\mathbf{r})], \quad (5)$$

which is composed of two-halo central-central (cen-cen) pairs

$$n_{\text{cc-pair}, 2\text{h}} = \frac{1}{2} \sum_{i \neq j} [\bar{n}_i \langle N_{\text{cen}}(M_i) \rangle] [\bar{n}_j \langle N_{\text{cen}}(M_j) \rangle d^3\mathbf{r}] \times [1 + \xi_{\text{hh}, \text{cc}}(\mathbf{r}; M_i, M_j)], \quad (6)$$

two-halo cen-sat pairs

$$n_{\text{cs-pair}, 2\text{h}} = \sum_{i \neq j} [\bar{n}_i \langle N_{\text{cen}}(M_i) \rangle] [\bar{n}_j \langle N_{\text{sat}}(M_j) \rangle d^3\mathbf{r}] \times [1 + \xi_{\text{hh}, \text{cs}}(\mathbf{r}; M_i, M_j)], \quad (7)$$

and two-halo sat-sat pairs

$$n_{\text{ss-pair}, 2\text{h}} = \frac{1}{2} \sum_{i \neq j} [\bar{n}_i \langle N_{\text{sat}}(M_i) \rangle] [\bar{n}_j \langle N_{\text{sat}}(M_j) \rangle d^3\mathbf{r}] \times [1 + \xi_{\text{hh}, \text{ss}}(\mathbf{r}; M_i, M_j)]. \quad (8)$$

In each of equations (6)–(8), the summation is over all halo mass bins (i.e.  $i = 1, \dots, N$  and  $j = 1, \dots, N$ ). The three correlation functions on the RHS have the following meanings –  $\xi_{\text{hh}, \text{cc}}(\mathbf{r}; M_i, M_j)$  is just the two-point cross-correlation function between ‘centres’ (positions to put central galaxies) of haloes of masses  $M_i$  and  $M_j$  (cen-cen);  $\xi_{\text{hh}, \text{cs}}(\mathbf{r}; M_i, M_j)$  is the two-point cross-correlation function between the ‘centres’ of  $M_i$  haloes and the satellite tracer particles in the (extended)  $M_j$  haloes (cen-sat);  $\xi_{\text{hh}, \text{ss}}(\mathbf{r}; M_i, M_j)$  is the two-point cross-correlation function between satellite tracer particles in the (extended)  $M_i$  haloes and those in the (extended)  $M_j$  haloes (sat-sat). With  $n_{\text{pair}, 2\text{h}} = n_{\text{cc-pair}, 2\text{h}} + n_{\text{cs-pair}, 2\text{h}} +$

$n_{\text{ss-pair}, 2\text{h}}$ , we reach the final expression for the two-halo term,

$$\xi_{\text{gg}}^{\text{2h}}(\mathbf{r}) = \sum_{i \neq j} \frac{\bar{n}_i \bar{n}_j}{\bar{n}_g^2} \langle N_{\text{cen}}(M_i) \rangle \langle N_{\text{cen}}(M_j) \rangle \xi_{\text{hh}, \text{cc}}(\mathbf{r}; M_i, M_j) + \sum_{i \neq j} 2 \frac{\bar{n}_i \bar{n}_j}{\bar{n}_g^2} \langle N_{\text{cen}}(M_i) \rangle \langle N_{\text{sat}}(M_j) \rangle \xi_{\text{hh}, \text{cs}}(\mathbf{r}; M_i, M_j) + \sum_{i \neq j} \frac{\bar{n}_i \bar{n}_j}{\bar{n}_g^2} \langle N_{\text{sat}}(M_i) \rangle \langle N_{\text{sat}}(M_j) \rangle \xi_{\text{hh}, \text{ss}}(\mathbf{r}; M_i, M_j). \quad (9)$$

Equations (1), (3), and (9) lead to the method we propose. The quantities related to galaxy occupancy are specified by the HOD/CLF parameterization one chooses, while those related to haloes are from the simulation, independent of the HOD/CLF parameterization. We therefore can prepare tables for  $\bar{n}_i$ ,  $f_{\text{cs}}(\mathbf{r}; M_i)$ ,  $f_{\text{ss}}(\mathbf{r}; M_i)$ ,  $\xi_{\text{hh}, \text{cc}}(\mathbf{r}; M_i, M_j)$ ,  $\xi_{\text{hh}, \text{cs}}(\mathbf{r}; M_i, M_j)$ , and  $\xi_{\text{hh}, \text{ss}}(\mathbf{r}; M_i, M_j)$ . For a given set of HOD/CLF parameters, the predictions of galaxy 2PCFs can be obtained from performing the weighted summation over the tables. The tables are only prepared once, and we can then change the galaxy occupation as needed to compute galaxy 2PCFs for different galaxy samples and different sets of HOD/CLF parameters, which is much more efficient than populating galaxies into haloes by selecting particles.

Since summation is used to replace integration in the method, we need to choose narrow halo mass bins ( $d \log M = 0.01$  is usually sufficient, as shown in Section 3). The  $\bar{n}_i$  table represents the halo mass function. To prepare the other tables that depend on pair separation, the bins of pair separation  $\mathbf{r}$  are best chosen to match the ones used in the measurements from observational data, which would naturally avoid any discrepancy related to the finite bin sizes. For haloes in each mass bin, the  $f_{\text{cs}}$  and  $f_{\text{ss}}$  tables can be computed by using either all the particles in the haloes with the specified distribution or a random subset. For  $\xi_{\text{hh}, \text{cc}}$ ,  $\xi_{\text{hh}, \text{cs}}$ , and  $\xi_{\text{hh}, \text{ss}}$ , we effectively compute the halo-halo two-point cross-correlation function with different definitions of halo positions. For  $\xi_{\text{hh}, \text{cc}}$ , halo positions are defined by our choice of ‘centres’. For  $\xi_{\text{hh}, \text{cs}}(\mathbf{r}; M_i, M_j)$ , we choose ‘centres’ for  $M_i$  haloes and positions of arbitrary tracer particles in  $M_j$  haloes. For  $\xi_{\text{hh}, \text{ss}}(\mathbf{r}; M_i, M_j)$ , positions of arbitrary tracer particles in both  $M_i$  and  $M_j$  haloes are chosen. We can use any number of tracer particles in each halo to do the calculation. For haloes with positions defined by the tracer particles, they can be thought as extended (with positions having a probability distribution). On large scales,  $\xi_{\text{hh}, \text{cc}}$ ,  $\xi_{\text{hh}, \text{cs}}$ , and  $\xi_{\text{hh}, \text{ss}}$  are the same, while on small scales,  $\xi_{\text{hh}, \text{cs}}$  and  $\xi_{\text{hh}, \text{ss}}$  are smoothed version of  $\xi_{\text{hh}, \text{cc}}$ . Note that in analytic models such differences are usually neglected. In computing the three halo-halo correlation functions, we do not need to construct random catalogs to find out the pair counts from a uniform distribution – in the volume  $V_{\text{sim}}$  of the simulation with periodic boundary conditions, the counts of cross-pairs at separation in the range  $\mathbf{r} \pm d\mathbf{r}/2$  between two randomly distributed populations with number densities  $\bar{n}_i$  and  $\bar{n}_j$  are simply  $(\bar{n}_i V_{\text{sim}})(\bar{n}_j d^3\mathbf{r})$ . Making use of this fact can greatly reduce the computational expense in preparing the tables.

For the redshift-space tables, in addition to the halo velocities, one needs to specify the velocity distribution of galaxies inside haloes, which can be different from that of dark matter particles (a.k.a. velocity bias; e.g. [Berlind & Weinberg 2002](#)). The difference can be parameterized by central and satellite velocity bias parameters (e.g. [Guo et al. 2015a](#)). For a set of central and satellite velocity bias parameters and with a choice of the line-of-sight direction, we can obtain the redshift-space positions of the central galaxy and satellite tracer particles according to halo velocities

and central and satellite galaxy velocity distributions inside haloes, and the redshift-space tables can be computed. We suggest to prepare tables for different sets of central and satellite velocity bias parameters and interpolate among tables to probe the velocity bias parameter space, as is done in [Guo et al. \(2015a\)](#).

Multipole moments of redshift-space galaxy 2PCFs are usually modelled. We can derive the corresponding tables by computing the corresponding multipole moments of  $f_{cs}$ ,  $f_{ss}$ ,  $\xi_{hh,cc}$ ,  $\xi_{hh,cs}$ , and  $\xi_{hh,ss}$ . In such a case,  $r$  is expressed by  $s = |\mathbf{r}|$  and  $\mu$ , the cosine of the angle between  $\mathbf{r}$  and the line-of-sight direction. In the integration (summation) for obtaining the multipoles, the bins of  $\mu$  match those used in observational measurements to remove any finite-bin-size effect.

For modelling the projected 2PCF  $w_p$ , a corresponding set of tables can be obtained by integrating the redshift-space tables over the line-of-sight separation. The integration is done in the same way as in the measurements with data to avoid any finite-bin-size effect, summing over the same line-of-sight bins (with the same bin size) up to the same maximum line-of-sight separation.

## 2.2 Case with Subhaloes

The SHAM method uses more information from (high-resolution) simulations, including both distinct haloes and subhaloes identified inside distinct haloes, where the distinct haloes refer to haloes that are not subhaloes of another halo. Distinct haloes are also referred to as haloes, main haloes, or host haloes. Central galaxies are hosted by distinct haloes at the centres, while satellite galaxies are in subhaloes. Before merging into distinct haloes, subhaloes are distinct haloes themselves. The SHAM method generally works in the following way. By adopting one property, subhaloes and distinct haloes can be treated as a unified entity. For distinct haloes, the property is evaluated at the time of interest. For subhaloes, it becomes common practice to evaluate the property at the time when subhaloes were still distinct haloes. The properties commonly used include mass ( $M_{acc}$ ) at the time a subhalo was accreted into a host halo, maximum circular velocity  $V_{acc}$  at the time of accretion, and peak maximum circular velocity  $V_{peak}$  over the history of the subhalo as a distinct halo. The connection between haloes/subhaloes and galaxies is established by rank ordering haloes/subhaloes according to the given property and galaxies according to one certain property (e.g. luminosity or stellar mass). When normalized to the same survey/simulation volume, halo/subhalo and galaxy of the same rank are linked. A more general treatment also accounts for the scatter between the halo/subhalo property and the galaxy property. The simple procedure of linking light (galaxies) to matter (haloes/subhaloes) can provide a reasonable interpretation of galaxy clustering trend and enable a study of galaxy evolution (e.g. [Conroy et al. 2006](#); [Conroy & Wechsler 2009](#); [Behroozi et al. 2013](#); [Reddick et al. 2013](#)).

We generalize the idea in Section 2.1 to the subhalo case, extending the SHAM model and making it efficient to model galaxy clustering. The model allows the scatter between the halo/subhalo property and the galaxy property to be different for distinct haloes (central galaxies) and subhaloes (satellite galaxies). We use mass as the halo/subhalo property variable here, which can be understood as the mass at accretion ( $M_{acc}$ ). However, it can be replaced by any property one chooses to adopt, e.g.  $V_{acc}$  and  $V_{peak}$ . A halo/subhalo method following a similar spirit of pair decomposition to model the projected galaxy 2PCF and weak lensing signal is presented in [Neistein et al. \(2011\)](#) and [Neistein & Khochfar \(2012\)](#).

Compared to the commonly used SHAM method that con-

nects the whole range of galaxy property and halo/subhalo property, the method presented here works for each individual galaxy sample. To some degree, it is formulated in an HOD/CLF-like form, with distinct haloes and subhaloes as tracers of central and satellite galaxies, respectively. It is no longer limited to *abundance* matching. Instead, the method can be used to fit both galaxy abundance and galaxy clustering (2PCFs).

For a given galaxy sample, the scatter between halo/subhalo property and galaxy property means that not all haloes/subhaloes are fully occupied by these galaxies, which can be characterized by the probability of occupancy (or the smaller-than-unity mean occupation number). Denote the mean occupation number of central galaxies in distinct haloes of mass  $M_h$  as  $p_{cen}(M_h)$  and that of satellite galaxies in subhaloes of mass  $M_s$  as  $p_{sat}(M_s)$ . The same bins of mass are adopted for  $M_h$  and  $M_s$ . In principle we do not need to differentiate  $M_s$  and  $M_h$ , since the scripts of ‘c’ (cen) and ‘s’ (sat) below make the situation self-explanatory. Let the mean number densities of distinct haloes and subhaloes in the mass bin  $\log M_i \pm d \log M_i/2$  be  $\bar{n}_{h,i}$  and  $\bar{n}_{s,i}$ , respectively.

For a given sample of galaxies, with a model of  $p_{cen}(M)$  and  $p_{sat}(M)$ , the mean number density of galaxies  $\bar{n}_g$  is computed as

$$\bar{n}_g = \sum_i [\bar{n}_{h,i} p_{cen}(M_i) + \bar{n}_{s,i} p_{sat}(M_i)]. \quad (10)$$

With a similar decomposition as in equation (9), the galaxy 2PCF can be computed as

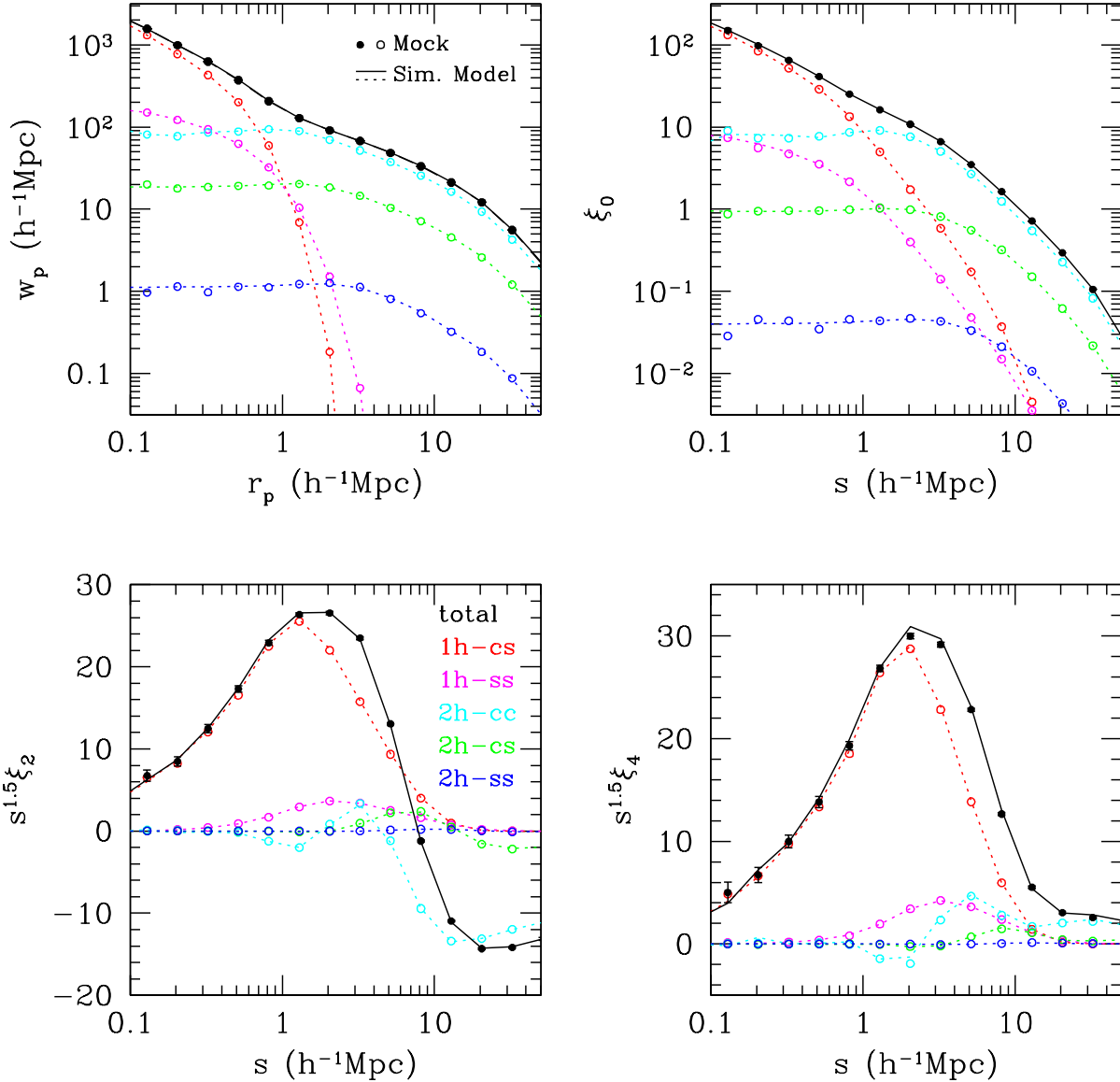
$$\begin{aligned} \xi_{gg}(\mathbf{r}) = & \sum_{i,j} \frac{\bar{n}_{h,i} \bar{n}_{h,j}}{\bar{n}_g^2} p_{cen}(M_i) p_{cen}(M_j) \xi_{hh}(\mathbf{r}; M_i, M_j) \\ & + \sum_{i,j} 2 \frac{\bar{n}_{h,i} \bar{n}_{s,j}}{\bar{n}_g^2} p_{cen}(M_i) p_{sat}(M_j) \xi_{hs}(\mathbf{r}; M_i, M_j) \\ & + \sum_{i,j} \frac{\bar{n}_{s,i} \bar{n}_{s,j}}{\bar{n}_g^2} p_{sat}(M_i) p_{sat}(M_j) \xi_{ss}(\mathbf{r}; M_i, M_j), \end{aligned} \quad (11)$$

which simply states that the total number of galaxy pairs is the sum of cen-cen, cen-sat, and sat-sat pairs. The three correlation functions on the RHS have the following meanings –  $\xi_{hh}(\mathbf{r}; M_i, M_j)$  is just the two-point cross-correlation function between centres of distinct haloes of masses  $M_i$  and  $M_j$ ;  $\xi_{hs}(\mathbf{r}; M_i, M_j)$  is the two-point cross-correlation function between centres of  $M_i$  distinct haloes and those of  $M_j$  subhaloes;  $\xi_{ss}(\mathbf{r}; M_i, M_j)$  is the two-point cross-correlation function between centres of subhaloes of masses  $M_i$  and  $M_j$ . Unlike the particle case in Section 2.1, there are no explicit one-halo and two-halo terms here (though they can be derived), and the  $i \neq j$  condition is not imposed in the summation.

The quantities  $p_{cen}(M)$  and  $p_{sat}(M)$  come from the occupation function model, which is up to our choice of parameterization for the sample of galaxies. In this halo/subhalo-based method, we only need to prepare tables for  $\bar{n}_{h,i}$ ,  $\bar{n}_{s,i}$ ,  $\xi_{hh}(\mathbf{r}; M_i, M_j)$ ,  $\xi_{hs}(\mathbf{r}; M_i, M_j)$ , and  $\xi_{ss}(\mathbf{r}; M_i, M_j)$ .

As with the tables using particles (Section 2.1), for redshift-space 2PCF or multipole moments, tables for different sets of central and satellite velocity bias parameters can be prepared. For each set, haloes and subhaloes are shifted to redshift-space positions for calculation. Tables can also be generated for modelling the projected 2PCF  $w_p$ . The procedures and bins used in the measurements should be followed so that the model and measurements are made fully consistent.



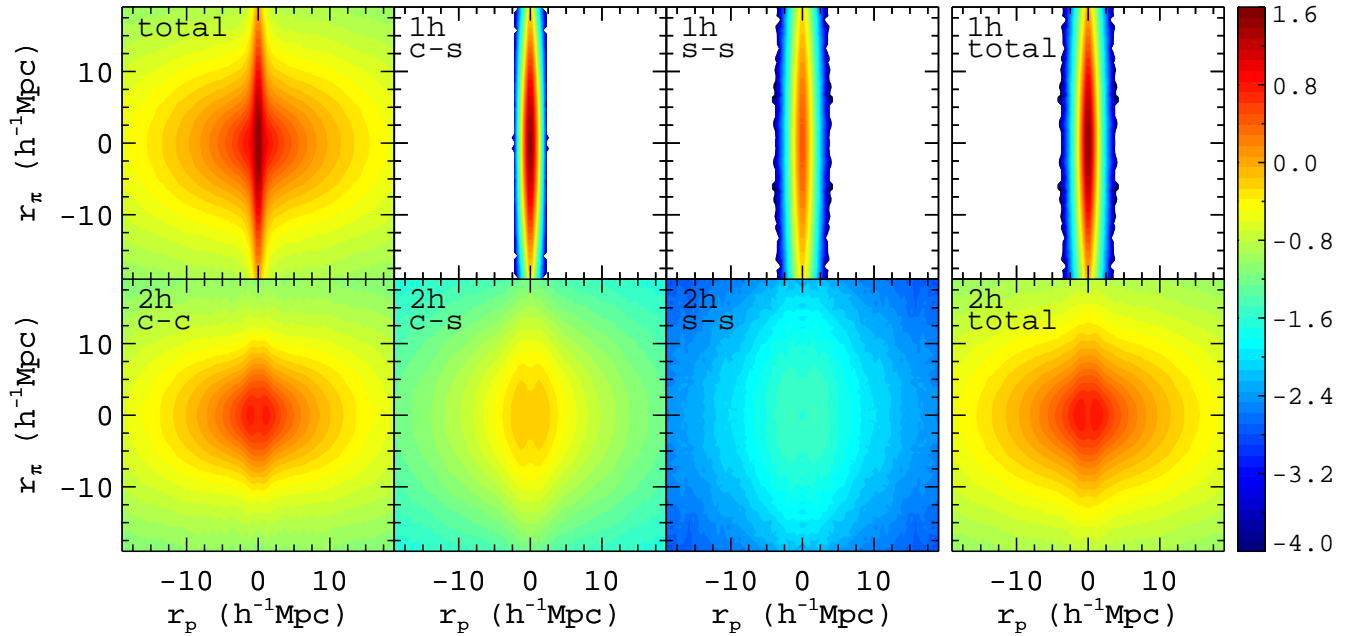


**Figure 1.** Decomposition of the projected galaxy 2PCF  $w_p$  and redshift-space 2PCF multipoles  $\xi_0$ ,  $\xi_2$ , and  $\xi_4$  into the various one-halo and two-halo components (one-halo cen-sat, one-halo sat-sat, two-halo cen-cen, two-halo cen-sat, and two-halo sat-sat). The circles are measurements from 100 mock galaxy catalogs constructed by populating galaxies into dark matter halos in the simulation, according to the set of fiducial HOD parameters. The curves are calculations with the method introduced in this paper. See text for more details.

### 3 AN EXAMPLE APPLICATION AND THE REDSHIFT-SPACE 2PCF DECOMPOSITION

The method developed here has been successfully applied to model projected and redshift-space 2PCFs of SDSS and SDSS-III galaxies on small to intermediate scales (e.g. Guo et al. 2015a,b,c) and to compare HOD and SHAM models (Guo et al. in prep.). As the method is built on the basis of decomposition of galaxy 2PCFs, here we provide an example to illustrate the different 2PCF components. In particular, we show the components for the redshift-space 3D 2PCF and the manifestation of redshift-space distortions in each component to have a better understanding of the redshift-space 2PCFs within the HOD framework. In addition, we also investigate how redshift-space 2PCFs help with HOD constraints, including the inference of the galaxy velocity distribution inside haloes.

The example adopts HOD parameters for the sample of  $z \sim 0.5$  CMASS galaxies in the the SDSS-III Baryon Oscillation Spectroscopic Survey (BOSS; Dawson et al. 2013). With spherical overdensity haloes and halo particles from the  $z = 0.53$  output of the MultiDark simulation (MDR1; Prada et al. 2012; Riebe et al. 2013), we create tables for halo properties, including halo number density  $\bar{n}$  (i.e. halo mass function), projected 2PCF  $w_p$ , redshift-space 2PCF monopole  $\xi_0$ , quadrupole  $\xi_2$ , and hexadecapole  $\xi_4$ . We choose the position of the potential minimum as the centre of each halo for putting the central galaxy and halo particles as tracers of satellites. Each of  $w_p$  and  $\xi_{0/2/4}$  has five components (one-halo cen-sat, one-halo sat-sat, two-halo cen-cen, two-halo cen-sat, and two-halo sat-sat). To generate the  $w_p(r_p)$  tables, we measure  $\xi(r_p, r_\pi)$  for each component and for each combination of halo mass bins and sum over the  $r_\pi$  direction, where  $r_p$  and  $r_\pi$  are the pair separations in the directions perpendicular and parallel to the



**Figure 2.** Decomposition of the 3D redshift-space 2PCF  $\xi(r_p, r_\pi)$  into the various one-halo and two-halo components (one-halo cen-sat, one-halo sat-sat, two-halo cen-cen, two-halo cen-sat, and two-halo sat-sat). The plot is based on the average measurements from 100 mock galaxy catalogs constructed by populating galaxies into dark matter halos in the simulation, according to the set of fiducial HOD parameters. The color scale shows  $\xi(r_p, r_\pi)$  in logarithmic scale. See text for more details.

line-of-sight direction (chosen to be one principle direction of the simulation box). To generate the  $\xi_{0/2/4}$  tables, we measure  $\xi(s, \mu)$  for each component and for each combination of halo mass bins and form the multipoles by integrating over  $\mu$ , where  $s$  is the redshift-space pair separation and  $\mu$  the cosine of the angle between pair displacement and the line-of-sight direction. Following the setup in the observational measurements (Guo et al. 2015a), we have 19 bins for  $r_p$  and  $s$  uniformly spaced in logarithmic space, 50 linearly spaced bins in  $r_\pi$  and 20 linearly spaced bins in  $\mu$ . For halo mass bins, we use  $d \log M = 0.01$ . We construct tables for 5 bins of central velocity bias parameter  $\alpha_c$  and 8 bins of satellite velocity bias parameter  $\alpha_s$ , respectively. The total size of the final set of tables is about 10GB. That is, the information in the high-resolution simulation output relevant for modelling projected and redshift-space 2PCFs of galaxies has been tremendously compressed, making the modelling tractable even with a desktop computer.

For the HOD, we adopt the common parameterization for a sample of galaxies above a luminosity threshold (Zheng et al. 2005, 2007). The mean occupation function of central galaxies in haloes of mass  $M$  is

$$\langle N_{\text{cen}}(M) \rangle = \frac{1}{2} \left[ 1 + \text{erf} \left( \frac{\log M - \log M_{\text{min}}}{\sigma_{\log M}} \right) \right], \quad (12)$$

where erf is the error function. For the mean occupation function of satellite galaxies, we use

$$\langle N_{\text{sat}}(M) \rangle = \langle N_{\text{cen}}(M) \rangle \left( \frac{M - M_0}{M'_1} \right)^\alpha. \quad (13)$$

The number of satellites in haloes of mass  $M$  is assumed to follow the Poisson distribution with the above mean. In addition, for modelling redshift-space 2PCFs, we have two additional HOD parameters  $\alpha_c$  and  $\alpha_s$  for central and satellite velocity bias. Essentially,  $\alpha_c$  ( $\alpha_s$ ) is the ratio of the velocity dispersion of central (satellite) galaxies to that of dark matter particles inside halos (see Guo et al.

2015a). For the fiducial model, we adopt the set of parameters that fit the projected and redshift-space 2PCFs for the CMASS sample in Guo et al. (2015a) –  $\log M_{\text{min}} = 13.36$ ,  $\sigma_{\log M} = 0.64$ ,  $\log M_0 = 13.20$ ,  $\log M'_1 = 14.23$ ,  $\alpha = 1.05$ ,  $\alpha_c = 0.30$ , and  $\alpha_s = 0.91$ . Halo masses are in units of  $h^{-1} M_\odot$ .

With the tables and the fiducial HOD parameters, we follow equations (1), (3), and (9) to compute all the components of  $w_p$  and  $\xi_{0/2/4}$ . For the purpose of a sanity check, we also measure the components from 100 mock galaxy catalogs. The mock catalogs are generated from populating haloes in the simulation by putting central galaxies at the potential minimum in haloes and drawing random dark matter particles as satellite galaxies, in accordance with the occupation distributions and velocities set by the fiducial HOD parameters. For the purpose of comparison with the model based on the tables, we decompose the galaxy 2PCF (either  $w_p$  or  $\xi_{0/2/4}$ ) measured in the mock catalogs into five components,

$$\begin{aligned} \xi_{\text{gg}}(\mathbf{r}) = & 2 \frac{\bar{n}_{\text{cs-pair}}}{\bar{n}_g^2} f_{\text{cs}}(\mathbf{r}) + 2 \frac{\bar{n}_{\text{ss-pair}}}{\bar{n}_g^2} f_{\text{ss}}(\mathbf{r}) \\ & + \frac{\bar{n}_c^2}{\bar{n}_g^2} \xi_{\text{cc}}(\mathbf{r}) + 2 \frac{\bar{n}_c \bar{n}_s}{\bar{n}_g^2} \xi_{\text{cs}}(\mathbf{r}) + \frac{\bar{n}_s^2}{\bar{n}_g^2} \xi_{\text{ss}}(\mathbf{r}). \end{aligned} \quad (14)$$

The first two terms on the RHS are one-halo terms –  $\bar{n}_{\text{cs-pair}}$  and  $\bar{n}_{\text{ss-pair}}$  are the mean number densities of one-halo cen-sat pairs and one-halo sat-sat pairs measured in the mock catalogs, and  $f_{\text{cs}}$  and  $f_{\text{ss}}$  are the normalized average distributions of one-halo cen-sat and sat-sat pairs in the mock. The last three terms on the RHS are two-halo terms –  $\bar{n}_c$  and  $\bar{n}_s$  are the mean number densities of central and satellite galaxies in the mock, and  $\xi_{\text{cc}}$ ,  $\xi_{\text{cs}}$ ,  $\xi_{\text{ss}}$  are the 2PCFs by counting only two-halo cen-cen, cen-sat, and sat-sat pairs (Zu et al. 2008).

Figure 1 shows the decomposition of  $w_p$  and  $\xi_{0/2/4}$  for the fiducial model. As expected, the calculations from the simulation-based method (curves) agree with the measurements from the mock

catalogs (circles), which is reassuring. For the projected 2PCF  $w_p$  (top-left panel), the one-halo cen-sat term (red) dominate the small-scale signal. The one-halo sat-sat term (magenta) extends to larger scales, since the maximum sat-sat pair separation in a halo is the diameter of the halo, twice that of the cen-sat pair separation. Owing to the low satellite fraction ( $f_{\text{sat}} \sim 7\%$ ) of this sample of galaxies, the contribution of the one-halo sat-sat pairs to  $w_p$  is overall small, but noticeable around  $1h^{-1}\text{Mpc}$ , the one-halo to two-halo term transition scales. On large scales, the three two-halo terms have a similar shape, since they essentially follow the halo-halo correlation. The flattening towards small scales are caused by the halo exclusion effect. Compared to the two-halo cen-cen component, the two-halo cen-sat is smoothed on small scales, since each halo contributing the satellite of the cen-sat pair on average is extended instead of a point source (the case for the halo contributing the central galaxy of the pair) as a result of the spatial distribution of satellites inside haloes. The two-halo sat-sat term is even more smoothed, since every halo becomes extended. To see the relative contribution of each term to the large-scale 2PCF, we note that in equation (14),  $\xi_{cc} \propto b_c^2$ ,  $\xi_{cs} \propto b_c b_s$ , and  $\xi_{ss} \propto b_s^2$  on large scales, where  $b_c$  and  $b_s$  are the large-scale bias factors for central and satellite galaxies, respectively. Since satellites on average reside in more massive haloes than central galaxies, the value of  $b_s$  is higher than that of  $b_c$  (roughly by tens of per cent for luminosity-threshold samples). From equation (14), we see that the relative contributions to the large-scale 2PCF from the two-halo cen-cen, cen-sat, and sat-sat terms are  $1 : 2f_n f_b : (f_n f_b)^2$ , with  $f_n = \bar{n}_s / \bar{n}_c = f_{\text{sat}} / (1 - f_{\text{sat}})$  the satellite to central galaxy number density ratio and  $f_b = b_s / b_c$  the satellite to central galaxy bias ratio. For the sample we consider, the ratios are  $1 : 25\% : 1.6\%$ . For lower luminosity samples with higher satellite fractions, we expect the contributions from the two-halo cen-sat and sat-sat to be substantially higher.

The decomposition of the redshift-space 2PCF monopole  $\xi_0$  (top-right panel) and the relative amplitudes of the various terms are similar to the case of  $w_p$ . The bottom two panels show the case of quadrupole  $\xi_2$  and hexadecapole  $\xi_4$ , and a factor  $s^{1.5}$  is multiplied for each term so that both the small-scale and large-scale signals can reasonably show up. The Fingers-of-God effect (Jackson 1972; Huchra 1988) from one-halo terms causes a positive quadrupole. In the  $\xi_2$  panel, we see that the influence of the one-halo terms can extend to about  $10h^{-1}\text{Mpc}$  in the quadrupole. The negative quadrupole on large scales manifests the Kaiser effect (Kaiser 1987; Hamilton 1992) caused by the coherent motion of haloes, falling into overdense regions and streaming out of underdense regions. The two-halo cen-cen term dominates the large-scale quadrupole, but the cen-sat term is also important. Both terms show low positive quadrupole signals toward small scales caused by the random motion of haloes (and galaxies). The two-halo sat-sat term makes an almost negligible contribution to the quadrupole on all scales. The hexadecapole  $\xi_4$  (bottom-right panel) are mostly positive from all components. The relative contributions from different components are similar to the quadrupole case.

The projected 2PCF and the redshift-space 2PCF multipoles are usually the quantities to model. The 3D redshift-space 2PCF measurements are commonly displayed as contours of  $\xi(r_p, r_\pi)$ , which make the redshift-space distortion effects on all scales easily visualized. It would be instructive to have the corresponding one-halo and two-halo components to gain a better intuition about the redshift-space distortions. Figure 2 shows such a decomposition measured from the mock catalogs, which can also be calculated using the  $\xi(r_p, r_\pi)$  component tables.

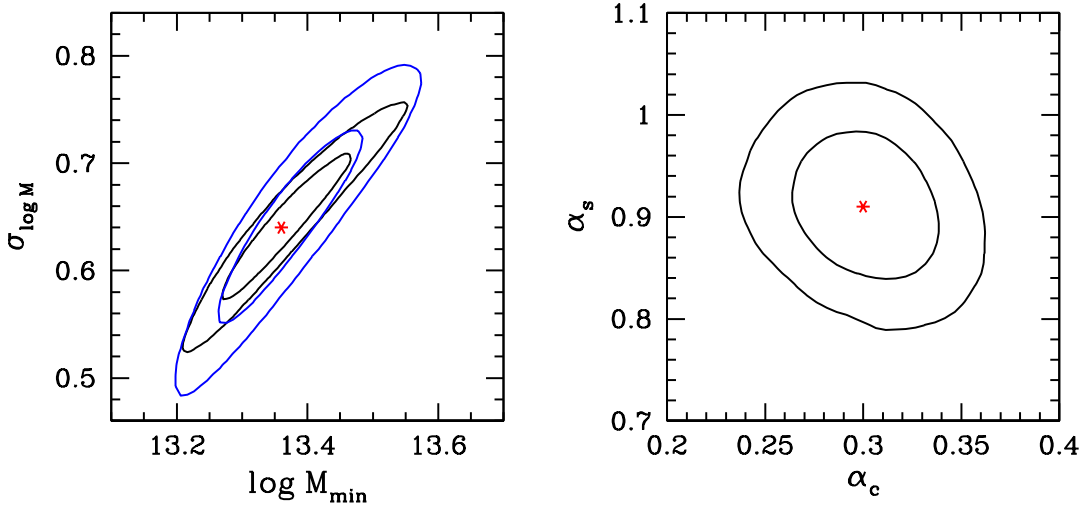
The leftmost panel shows the total redshift-space 2PCF of the

sample, with the Fingers-of-God and Kaiser effects clearly seen. The Fingers-of-God effect, limited to small transverse separation  $r_p$ , is mainly contributed by the one-halo terms (two middle panels on the top). The one-halo sat-sat component appears to be more extended than the one-halo cen-sat component in both the transverse and the line-of-sight direction. In the transverse direction, it can be explained by the fact that the largest one-halo sat-sat (cen-sat) pair separation is about the diameter (radius) of the largest haloes. In the line-of-sight direction, the elongation is mainly a result of galaxy motion inside haloes. The relative line-of-sight velocity of sat-sat pairs are higher than that of cen-sat pairs, causing the one-halo sat-sat component to be more extended (shallower profile as a function of  $r_\pi$ ). The total one-halo term (rightmost panel on the top) is dominated by the cen-sat and sat-sat component at small  $r_p$  and slightly large  $r_p$ , respectively.

The three two-halo components and the total two-halo term are shown in the bottom panels of Figure 2. In each component, the double-hump feature at small  $r_p$  reflects the halo-exclusion effect. The effect would lead to a hole at the centre if the real-space 2PCF were plotted here. The shift in the line-of-sight galaxy positions in redshift space from galaxy peculiar motion makes the hole partially filled. The two-halo cen-cen component shows an overall Kaiser squashing effect along the line of sight. However, the contours at small  $r_p$  are elongated along the line of sight, like the Fingers-of-God effect. This is caused by the random motion of haloes and that of central galaxies with respect to haloes (i.e. a non-zero central velocity bias). The two-halo cen-sat component shows a much stronger line-of-sight elongation up to a few Mpc in  $r_p$ . The reason lies in the motion of satellites inside haloes, which causes the average redshift-space distribution of satellites appears extended along the line of sight in an average halo hosting the satellites of the two-halo cen-sat pairs. The line-of-sight elongation pattern is even stronger in the two-halo sat-sat component – the correlation of elongated haloes (as a result of the redshift-space spatial distribution of satellites inside haloes) completely suppresses the Kaiser effect even on the largest scales shown here ( $\sim 20h^{-1}\text{Mpc}$ ). The total two-halo term is dominated by the cen-cen component with a substantial contribution from the cen-sat component. The sat-sat component does not make an important contribution for this sample. As discussed before, we expect the two-halo cen-sat and sat-sat components to become more important for galaxy samples with lower luminosity thresholds and higher satellite fractions.

Overall, for the 3D redshift 2PCF  $\xi(r_p, r_\pi)$  different components of the one-halo and two-halo terms have different transverse range of the line-of-sight elongation. The profile along the line of sight also depends on the type of pairs in consideration, becoming increasingly shallower from cen-cen, cen-sat, to sat-sat components. For each component, the streaming model (e.g. Peebles 1980) usually adopted in simple models of redshift-space distortions should work well, which is kind of a convolution of the real-space 2PCF with a velocity dispersion kernel. For the total redshift-space 2PCF, our results indicate that it is hard to use a single velocity dispersion kernel to accurately model the redshift-space distortion effect. The different components are needed if one wishes to develop an accurate analytic model (e.g. Tinker 2007).

Finally, we investigate the constraints on the HOD parameters from projected and redshift-space 2PCFs. The 2PCFs predicted from the fiducial set of HOD parameters are used as the input measurements, and the full covariance matrix from Guo et al. (2015a) measured from the CMASS data is adopted. The model uncertainty caused by the finite volume of the simulation is also accounted for by rescaling the covariance matrix (see Appendix A). We em-



**Figure 3.** *Left:* Constraints on  $\log M_{\min}$  and  $\sigma_{\log M}$  from the 2PCFs with the fiducial galaxy sample. The model 2PCFs are calculated with method introduced in this paper. Blue and black contours are for the cases of modelling  $w_p$  only and jointly modelling  $w_p + \xi_{0/2/4}$ , respectively. The 68.3% and 95.4% confidence levels are shown for each case. *Right:* Constraints on the central and satellite velocity bias parameters ( $\alpha_c$  and  $\alpha_s$ ) for the fiducial galaxy sample from jointly modelling  $w_p + \xi_{0/2/4}$ . The red asterisk in each panel indicates the value from the fiducial model.

ploy a Monte Carlo Markov Chain method to explore the parameter space of the 7 HOD parameters,  $M_{\min}$ ,  $\sigma_{\log M}$ ,  $M_0$ ,  $M'_1$ ,  $\alpha$ ,  $\alpha_c$ , and  $\alpha_s$ . We first model the projected 2PCF  $w_p$  only. The first five parameters related to the galaxy mean occupation function can be constrained, while there are virtually no constraints on the velocity bias parameters ( $\alpha_c$  and  $\alpha_s$ ) as the line-of-sight information is lost. We then jointly model  $w_p$  and the redshift-space 2PCF multipoles  $\xi_{0/2/4}$ . We find that redshift-space 2PCFs help tighten the constraints mainly in  $M_{\min}$  and  $\sigma_{\log M}$ , the two parameters for the mean occupation function of central galaxies. In the left panel of Figure 3, we compare the constraints (marginalized  $1\sigma$  and  $2\sigma$  contours) from  $w_p$  only (blue) and  $w_p + \xi_{0/2/4}$  (black). The constraints on the parameters for the mean occupation function of satellite galaxies are only slightly improved, mainly in  $M_0$ . In general, compared to the  $w_p$ -only case, redshift-space 2PCFs do not lead to a substantial improvement in the HOD parameters related to the occupation function. The reason may be related to the fact that the projected 2PCF  $w_p$  is not independent of the redshift-space 2PCFs, and that the information content in  $\xi_{0/2/4}$  to constrain the occupation-related parameters is largely overlapped with that in  $w_p$ . The correlated information in  $w_p$  and  $\xi_{0/2/4}$  is embedded in the covariance matrix. Therefore, when jointly modelling  $w_p$  and  $\xi_{0/2/4}$ , it is important to use the full covariance matrix including the covariances between  $w_p$  and  $\xi_{0/2/4}$  to avoid double counting the information content and artificially tightening the HOD constraints.

The redshift-space distortions are caused by the peculiar motion of galaxies. The peculiar motion of haloes is in the simulation and built in the tables. So modelling redshift-space 2PCFs lead to constraints of galaxy motion inside haloes, i.e. the central and satellite velocity bias parameters. The right panel of Figure 3 shows that velocity bias parameters can be clearly detected for the fiducial sample. Velocity bias parameters have been constrained from redshift-space clustering for the  $z \sim 0.5$  BOSS CMASS galaxies (Guo et al. 2015a,b; Reid et al. 2014) and  $z \sim 0.1$  SDSS Main galaxies (see Guo et al. 2015c and Guo et al. 2015d for applying the modelling method based on simulation particles and subhaloes,

respectively). More discussions on the velocity bias constraints and the implications can be found in Guo et al. (2015a).

#### 4 SUMMARY AND DISCUSSION

In this paper, we introduce a simulation-based method to accurately and efficiently model galaxy 2PCFs in projected and redshift spaces. The basic idea is to make use of a high-resolution simulation and tabulate all the halo information necessary for galaxy clustering calculation. Then on top of the tables, galaxy 2PCFs can be computed with the galaxy-halo relation specified by the HOD or CLF model. We also provide a version that applies to and extends the SHAM method. Based on the method, we also study the decomposition of the projected and redshift-space galaxy 2PCFs into different components according to the type of galaxy pairs.

The proposed method is accurate, since it is directly based on high-resolution simulations. The effects like halo exclusion, non-linear evolution, scale-dependent halo bias, and non-sphericity of haloes, which are difficult to deal with in analytic methods of computing galaxy 2PCFs, are all automatically accounted for in the simulation-based method. The method also breaks the 2PCFs into all the one-halo and two-halo components based on the nature of galaxy pairs and computes each component accurately, which are usually not the case in analytic methods (especially for the two-halo term). When building the tables, the same binning scheme (in pair separation and in angle) and the same integration procedure as used in the observation measurements are adopted, so there is no binning-related issue when comparing the model prediction with the measurements. The method is equivalent to measure the model galaxy 2PCFs from mock catalogs and is as accurate as what the mean mock catalog can achieve. The mock catalogs are constructed by populating galaxies (using tracer particles) to haloes identified in the simulation, according to the halo occupation specified by the HOD/CLF model. However, the method is more efficient, as it avoids the construction of mock catalogs and the measurement of the 2PCFs from the mocks. Instead, ‘populating galaxies’ and ‘measuring the 2PCFs’ are performed analytically within



the HOD/CLF framework. This greatly reduces the computational time and make it possible to efficiently explore the parameter space when modelling the 2PCF data.

A similar method working in Fourier space can be easily developed to model galaxy redshift-space power spectrum. The method can also be generalized to other clustering statistics, e.g. angular 2PCF of galaxies, two-point cross-correlation function of galaxies, and galaxy-galaxy lensing. Generalizing the method to three-point correlation function (3PCF) of galaxies is also possible. In principle, there are more components for the 3PCF – cen-sat-sat and sat-sat-sat triplets for the one-halo term, cen-(cen-sat), cen-(sat-sat), sat-(cen-sat), and sat-(sat-sat) triplets for the two-halo term (the pair in the parentheses is in the same halo), and cen-cen-cen, cen-cen-sat, cen-sat-sat, and sat-sat-sat triplets for the three-halo term. More importantly, compared to the 2PCF case, the dimension of each 3PCF component table will increase (e.g. two sides and the angle in between for a triangle configuration and three halo mass indices). To make such a method suitable for the 3PCF modelling, further simplification is necessary, e.g. through multipole or Fourier expansion (e.g. Szapudi 2004; Zheng 2004b; Slepian & Eisenstein 2015).

To make use of the high precision of small- to intermediate-scale 2PCFs measurements to help constrain cosmological parameters (e.g. Zheng & Weinberg 2007; Reid et al. 2014), a set of tables need to be prepared based on simulations with different cosmological parameters or by rescaling one simulation to different cosmological models (e.g. Zheng et al. 2002; Tinker et al. 2006; Angulo & White 2010; Reid et al. 2014; Guo et al. 2015c). Even with one cosmological model, there may be situations that need more tables. For example, in the particle-based model, random particles are selected to trace satellite galaxies by default. However, the difference between the spatial distributions of satellites and dark matter can be an additional parameter to be constrained. For such a purpose, one needs to build different sets of tables using tracer particles of different distributions. In either of the above cases (or any case that needs to extend the tables), the total size of the tables would have an order-of-magnitude increase. Compared with methods of directly populating simulations, such an increase in table size is still reasonable and manageable.

With one simulation, we do not have the global or ensemble average properties of haloes. That is, the model with one simulation has uncertainty caused by the finite volume effect. One can use multiple simulations with different realizations of the initial conditions to build the average tables, which reduces the model uncertainty. The model uncertainty should be included in modelling data. In Appendix A, we show that this can be done by rescaling the covariance matrix of the measurements based on the ratio of simulation and survey volume. For any simulation, the fluctuation modes with wavelengths longer than the box size are missing, so the application of our modelling method should be limited to scales much smaller than the simulation box size. This is particularly true for redshift-space distortion modelling, since the velocity field is more sensitive to large-scale modes than the density field.

In presenting the method, the halo variable is adopted to be halo mass (or characteristic velocity for the subhalo case) to build the tables. The corresponding HOD/CLF model assumes that the statistical properties of galaxies inside haloes only depend on halo mass, not on halo environment or growth history. Clustering of haloes at fixed mass is found to depend on the assembly history (a.k.a. assembly bias; e.g. Gao et al. 2005; Wechsler et al. 2006; Zhu et al. 2006; Jing et al. 2007). There is room for the galaxy content in haloes of fixed mass to depend on halo formation history,

which would affect galaxy clustering and HOD constraints (e.g. Zentner et al. 2014), although no clear evidence is found in hydrodynamic galaxy formation simulations (e.g. Berlind et al. 2003) or galaxy clustering measurements (e.g. Lin et al. 2015). As mentioned in Section 2, the halo variable in our method is not necessarily the halo mass. It can certainly be a set of variables, like halo mass plus a variable characterizing halo formation history (e.g. halo concentration or formation redshift). With tables built in terms of the set of variables, along with an HOD/CLF model depending on these variables, the simulation-based method works in the same way as presented in this paper. However, the efficiency of the method drops sharply when including more halo variables. The limitation is mainly set by the computation of the two-halo terms, where both the table size and computational time scale as  $\mathcal{O}(N^2)$ , with  $N$  the total number of bins in halo properties (e.g. with  $N_1$  halo mass bins and  $N_2$  halo formation time bins,  $N = N_1 N_2$ ). In practice, we may be barely able to accommodate the case of two halo variables, by choosing bin sizes to minimize the table size and computational cost without sacrificing the accuracy of the method. Before resorting to directly populating the simulations, a possible way of circumventing the limitation is to use some combination of halo variables, reducing the problem to one effective halo variable. Certainly further investigations are needed to find the appropriate combination(s).

A different approach to model galaxy clustering is through an emulator (e.g. Kwan et al. 2015). With this approach, galaxy correlation functions are first obtained with mock catalogs from  $N$ -body simulations, spanning a range of HOD parameters. Then the emulator works by interpolation to predict the galaxy correlation function for any given set of HOD parameters. Compared to the method we propose in this paper, the emulator can be extremely fast, since it only performs interpolations and avoids any calculation at the level of dark matter haloes. In principle, the emulator can be generalized to interpolate among the one-halo and two-halo component contributions to the 2PCFs. However, by construction, the emulator only operates with a certain HOD form and within a certain range of HOD parameters for the interpolation to work and for the accuracy to be under control. The method we propose performs direct calculations with clear physical meanings based on halo properties, and therefore it does not suffer from the above restrictions of an emulator.

With increasingly more precise measurements of galaxy clustering from forthcoming large galaxy surveys, such as DESI (Levi et al. 2013) and Euclid (Laureijs et al. 2011), we expect that the accurate and efficient modelling method introduced in this work and its generalizations will have great potentials and wide applications.

## ACKNOWLEDGMENTS

ZZ is partially supported by NSF grant AST-1208891. HG acknowledges the support of NSFC-11543003 and the 100 Talents Program of the Chinese Academy of Sciences.

## References

- Angulo R. E., White S. D. M., 2010, MNRAS, 405, 143
- Behroozi P. S., Wechsler R. H., Conroy C., 2013, ApJ, 770, 57
- Berlind A. A., Weinberg D. H., 2002, ApJ, 575, 587
- Berlind A. A. et al., 2003, ApJ, 593, 1

Blake C., Kazin E. A., Beutler F. et al., 2011, *MNRAS*, 418, 1707  
 Cacciato M., van den Bosch F. C., More S., Mo H., Yang X., 2013, *MNRAS*, 430, 767  
 Colless M., 1999, *RSPTA*, 357, 105  
 Conroy C., Wechsler R. H., Kravtsov A. V., 2006, *ApJ*, 647, 201  
 Conroy C., Wechsler R. H., 2009, *ApJ*, 696, 620  
 Dawson K. S. et al., 2013, *AJ*, 145, 10  
 Eisenstein D. J., Weinberg D. H., Agol E. et al. 2011, *AJ*, 142, 72  
 Feldman H. A., Kaiser N., Peacock J. A., 1994, *ApJ*, 426, 23  
 Gao L., Springel V., White S. D. M., 2005, *MNRAS*, 363, L66  
 Guo H. et al., 2015a, *MNRAS*, 446, 578  
 Guo H. et al., 2015b, *MNRAS*, 449, L95  
 Guo H. et al., 2015c, *MNRAS*, 453, 4368  
 Guo H. et al., 2015d, *arXiv:1508.07012*  
 Guo H. et al., 2014, *MNRAS*, 441, 2398  
 Hamilton A. J. S., 1992, *ApJL*, 385, L5  
 Huchra J. P., 1988, in Dickey J. M., ed., *Astronomical Society of the Pacific Conference Series Vol. 5, The Minnesota lectures on Clusters of Galaxies and Large-Scale Structure*. pp 41–70  
 Jackson J. C., 1972, *MNRAS*, 156, 1P  
 Jenkins A., Frenk C. S., White S. D. M. et al., 2001, *MNRAS*, 321, 372  
 Jing Y. P., Mo H. J., Boerner G., 1998, *ApJ*, 494, 1  
 Jing Y. P., Suto Y., Mo H. J., 2007, *ApJ*, 657, 664  
 Kaiser N., 1987, *MNRAS*, 227, 1  
 Kwan J., Heitmann K., Habib S., et al., 2015, *ApJ*, 810, 35  
 Laureijs R., Amiaux J., Arduini S. et al., 2011, *ArXiv e-prints*, [arXiv:1110.3193](https://arxiv.org/abs/1110.3193)  
 Levi M., Bebek C., Beers T. et al., 2013, *ArXiv e-prints*, [arXiv:1308.0847](https://arxiv.org/abs/1308.0847)  
 Lin Y.-T., Mandelbaum R., Huang Y.-H., Huang H.-J., Dalal N., Diemer B., Jian H.-Y., Kravtsov A., 2015, *ArXiv e-prints*, [arXiv:1504.07632](https://arxiv.org/abs/1504.07632)  
 Mo H. J., Jing Y. P., White S. D. M., 1996, *MNRAS*, 282, 1096  
 Nagai D., Kravtsov A. V., 2005, *ApJ*, 618, 557  
 Neistein E., Li C., Khochfar S. et al., 2011, *MNRAS*, 416, 1486  
 Neistein E., Khochfar S., 2012, *ArXiv e-prints*, [arXiv:1209.0463](https://arxiv.org/abs/1209.0463)  
 Parejko J. K. et al., 2013, *MNRAS*, 429, 98  
 Peacock J. A., Smith R. E., 2000, *MNRAS*, 318, 1144  
 Peebles P. J. E., 1980, *The large-scale structure of the universe*  
 Prada F., Klypin A. A., Cuesta A. J., Betancort-Rijo J. E., Primack J., 2012, *MNRAS*, 423, 3018  
 Press W. H., Schechter P., 1974, *ApJ*, 187, 425  
 Reddick R. M., Wechsler R. H., Tinker J. L., Behroozi P. S., 2013, *ApJ*, 771, 30  
 Reid B. A., Seo H.-J., Leauthaud A., Tinker J. L., White M., 2014, *MNRAS*, 444, 476  
 Reid B. A., White M., 2011, *MNRAS*, 417, 1913  
 Riebe K. et al., 2013, *AN*, 334, 691  
 Scoccimarro R., Sheth R. K., Hui L., Jain B., 2001, *ApJ*, 546, 20  
 Seljak U., 2000, *MNRAS*, 318, 203  
 Sheth R. K., Tormen G., 1999, *MNRAS*, 308, 119  
 Skibba R. A., Coil A. L., Mendez A. J., et al., 2015, *ApJ*, 807, 152  
 Slepian Z., Eisenstein D. J., 2015, *MNRAS*, 454, 4142  
 Smith R. E. et al., 2003, *MNRAS*, 341, 1311  
 Szapudi I., 2004, *ApJL*, 605, L89  
 Tegmark M., 1997, *Physical Review Letters*, 79, 3806  
 Tinker J. L., Weinberg D. H., Zheng Z., 2006, *MNRAS*, 368, 85  
 Tinker J. L., 2007, *MNRAS*, 374, 477  
 Tinker J. L., Weinberg D. H., Zheng Z., Zehavi I., 2005, *ApJ*, 631, 41  
 Tinker J., Kravtsov A. V., Klypin A. et al., 2008, *ApJ*, 688, 709

Tinker J. L., Robertson B. E., Kravtsov A. V. et al., 2010, *ApJ*, 724, 878  
 van den Bosch F. C., Yang X. H., Mo H. J., 2003, *MNRAS*, 340, 771  
 van den Bosch F. C., Mo H. J., Yang X. H., 2003, *MNRAS*, 345, 923  
 van den Bosch F. C., More S., Cacciato M., Mo H., Yang X., 2013, *MNRAS*, 430, 725  
 Wechsler R. H., Zentner A. R., Bullock J. S., Kravtsov A. V., Allgood B., 2006, *ApJ*, 652, 71  
 White M. et al., 2011, *ApJ*, 728, 126  
 Yang X., Mo H. J., van den Bosch F. C., 2003, *MNRAS*, 339, 1057  
 York D. G. et al., 2000, *AJ*, 120, 1579  
 Zehavi I. et al., 2005, *ApJ*, 621, 22  
 Zehavi I. et al., 2011, *ApJ*, 736, 59  
 Zentner A. R., Hearin A. P., van den Bosch F. C., 2014, *MNRAS*, 443, 3044  
 Zheng Z., Tinker J. L., Weinberg D. H., Berlind A. A., 2002, *ApJ*, 575, 617  
 Zheng Z., 2004a, *ApJ*, 610, 61  
 Zheng Z., 2004b, *ApJ*, 614, 527  
 Zheng Z. et al., 2005, *ApJ*, 633, 791  
 Zheng Z., Coil A. L., Zehavi I., 2007, *ApJ*, 667, 760  
 Zheng Z., Weinberg D. H., 2007, *ApJ*, 659, 1  
 Zheng Z., Zehavi I., Eisenstein D. J., Weinberg D. H., Jing Y. P., 2009, *ApJ*, 707, 554  
 Zhu G., Zheng Z., Lin W. P., Jing Y. P., Kang X., Gao L., 2006, *ApJL*, 639, L5  
 Zu Y., Weinberg D. H., 2013, *MNRAS*, 431, 3319  
 Zu Y., Zheng Z., Zhu G., Jing Y. P., 2008, *ApJ*, 686, 41

## APPENDIX A: COVARIANCE MATRIX WITH MODEL UNCERTAINTY

Let us consider the case that we use a model built on one simulation in a volume  $V_m$  ('m' for model) to interpret the observation obtained from a survey volume  $V_o$  ('o' for observation). What covariance matrix should we use to model the data? The covariance matrix estimated for the observation tells us the covariance in the observational data. However, the model is based on a simulation with a finite volume, and therefore it is not the global model or the model from ensemble average. The model itself has uncertainty, and the modelling needs to account for this. To derive the effective covariance matrix  $C^{\text{eff}}$  to be used in the modelling, let us define the  $i$ -th data point measured in the observational volume  $V_o$  as  $F_{o,i}^{V_o}$ , the  $i$ -th data point from the model with simulation volume  $V_m$  as  $F_{m,i}^{V_m}$ , and the global averages (or the ensemble averages) of the observational and model data points as  $F_{o,i}$  and  $F_{m,i}$ , respectively. Note that for an accurate model that reflects the reality, we have  $F_{m,i} = F_{o,i}$ . That is, the global model reproduces the global average observation.

The effective covariance matrix with model uncertainty included is then

$$C_{ij}^{\text{eff}} = \langle (F_{o,i}^{V_o} - F_{m,i}^{V_m}) (F_{o,j}^{V_o} - F_{m,j}^{V_m}) \rangle \quad (\text{A1})$$

$$= \langle [(F_{o,i}^{V_o} - F_{o,i}) - (F_{m,i}^{V_m} - F_{m,i})] [(F_{o,j}^{V_o} - F_{o,j}) - (F_{m,j}^{V_m} - F_{m,j})] \rangle \quad (\text{A2})$$

$$= \langle (F_{o,i}^{V_o} - F_{o,i}) (F_{o,j}^{V_o} - F_{o,j}) \rangle$$

$$\begin{aligned}
& + \langle (F_{m,i}^{V_m} - F_{m,i}) (F_{m,j}^{V_m} - F_{m,j}) \rangle \\
& + \langle (F_{o,i}^{V_o} - F_{o,i}) (F_{m,j}^{V_m} - F_{m,j}) \rangle \\
& + \langle (F_{m,i}^{V_m} - F_{m,i}) (F_{o,j}^{V_o} - F_{o,j}) \rangle.
\end{aligned} \tag{A3}$$

The symbol  $\langle \rangle$  denotes global/ensemble average over observations in volumes of  $V_o$  and over models in volumes of  $V_m$ . From (A1) to (A2), we make use of the above  $F_{m,i} = F_{o,i}$  relation. In (A3), the first term is the element  $C_{ij}^{V_o}$  of the covariance matrix for the measurements in volume  $V_o$ , the second term is the element  $C_{ij}^{V_m}$  of the covariance matrix for the measurements in volume  $V_m$  (since the model values can be regarded as mock measurements), and both the third and fourth terms are zero (since there is no correlation between observation measurements and mock measurements). We then have

$$\mathbf{C}^{\text{eff}} = \mathbf{C}^{V_o} + \mathbf{C}^{V_m}, \tag{A4}$$

and the result is expected and intuitive.

For power spectrum or 2PCF measurements, the covariance matrix element is inversely proportional to the volume (Feldman et al. 1994; Tegmark 1997). We can express the effective covariance matrix in equation (A4) in terms of the one estimated for the observation and the relative volume of the simulation and observation,

$$\mathbf{C}^{\text{eff}} = \left(1 + \frac{V_o}{V_m}\right) \mathbf{C}^{V_o}. \tag{A5}$$